

# Reporting guidelines for clinical trial reports for interventions involving artificial intelligence

The SPIRIT-AI and CONSORT-AI Working Group

DOI:

[10.1038/s41591-020-1034-x](https://doi.org/10.1038/s41591-020-1034-x)

License:

Creative Commons: Attribution (CC BY)

*Document Version*

Publisher's PDF, also known as Version of record

*Citation for published version (Harvard):*

The SPIRIT-AI and CONSORT-AI Working Group 2020, 'Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI extension', *Nature Medicine*, vol. 26, no. 9, pp. 1364–1374. <https://doi.org/10.1038/s41591-020-1034-x>

[Link to publication on Research at Birmingham portal](#)

## General rights

Unless a licence is specified above, all rights (including copyright and moral rights) in this document are retained by the authors and/or the copyright holders. The express permission of the copyright holder must be obtained for any use of this material other than for purposes permitted by law.

- Users may freely distribute the URL that is used to identify this publication.
- Users may download and/or print one copy of the publication from the University of Birmingham research portal for the purpose of private study or non-commercial research.
- User may use extracts from the document in line with the concept of 'fair dealing' under the Copyright, Designs and Patents Act 1988 (?)
- Users may not further distribute the material nor use it for the purposes of commercial gain.

Where a licence is displayed above, please note the terms and conditions of the licence govern your use of this document.

When citing, please reference the published version.

## Take down policy

While the University of Birmingham exercises care and attention in making items available there are rare occasions when an item has been uploaded in error or has been deemed to be commercially or otherwise sensitive.

If you believe that this is the case for this document, please contact [UBIRA@lists.bham.ac.uk](mailto:UBIRA@lists.bham.ac.uk) providing details and we will remove access to the work immediately and investigate.



OPEN

# Reporting guidelines for clinical trial reports for interventions involving artificial intelligence: the CONSORT-AI extension

Xiaoxuan Liu<sup>1,2,3,4,5</sup>, Samantha Cruz Rivera<sup>5,6,7</sup>, David Moher<sup>8,9</sup>, Melanie J. Calvert<sup>10,4,5,6,7,10,11,12</sup>, Alastair K. Denniston<sup>12,3,4,5,6,13</sup> ✉ and The SPIRIT-AI and CONSORT-AI Working Group\*

**The CONSORT 2010 statement provides minimum guidelines for reporting randomized trials. Its widespread use has been instrumental in ensuring transparency in the evaluation of new interventions. More recently, there has been a growing recognition that interventions involving artificial intelligence (AI) need to undergo rigorous, prospective evaluation to demonstrate impact on health outcomes. The CONSORT-AI (Consolidated Standards of Reporting Trials–Artificial Intelligence) extension is a new reporting guideline for clinical trials evaluating interventions with an AI component. It was developed in parallel with its companion statement for clinical trial protocols: SPIRIT-AI (Standard Protocol Items: Recommendations for Interventional Trials–Artificial Intelligence). Both guidelines were developed through a staged consensus process involving literature review and expert consultation to generate 29 candidate items, which were assessed by an international multi-stakeholder group in a two-stage Delphi survey (103 stakeholders), agreed upon in a two-day consensus meeting (31 stakeholders) and refined through a checklist pilot (34 participants). The CONSORT-AI extension includes 14 new items that were considered sufficiently important for AI interventions that they should be routinely reported in addition to the core CONSORT 2010 items. CONSORT-AI recommends that investigators provide clear descriptions of the AI intervention, including instructions and skills required for use, the setting in which the AI intervention is integrated, the handling of inputs and outputs of the AI intervention, the human–AI interaction and provision of an analysis of error cases. CONSORT-AI will help promote transparency and completeness in reporting clinical trials for AI interventions. It will assist editors and peer reviewers, as well as the general readership, to understand, interpret and critically appraise the quality of clinical trial design and risk of bias in the reported outcomes.**

Randomized controlled trials (RCTs) are considered the gold-standard experimental design for providing evidence of the safety and efficacy of an intervention<sup>1,2</sup>. Trial results, if adequately reported, have the potential to inform regulatory decisions, clinical guidelines and health policy. It is therefore crucial that RCTs are reported with transparency and completeness so that readers can critically appraise the trial methods and findings and assess the presence of bias in the results<sup>3–5</sup>.

The CONSORT statement provides evidence-based recommendations to improve the completeness of the reporting of RCTs. The statement was first introduced in 1996 and has since been widely endorsed by medical journals internationally<sup>5</sup>. Over the past two decades, it has undergone two updates and has demonstrated a substantial positive impact on the quality of RCT reports<sup>6,7</sup>. The most recent CONSORT 2010 statement provides a 25-item checklist of the minimum reporting content applicable to all RCTs, but it recognizes that certain interventions may require extension or elaboration of these items. Several such extensions exist<sup>8–13</sup>.

AI is an area of enormous interest with strong drivers to accelerate new interventions through to publication, implementation and market<sup>14</sup>. While AI systems have been researched for some time, recent advances in deep learning and neural networks have gained considerable interest for their potential in health applications. Examples of such applications are wide ranging and include AI systems for screening and triage<sup>15,16</sup>, diagnosis<sup>17–20</sup>, prognostication<sup>21,22</sup>, decision support<sup>23</sup> and treatment recommendation<sup>24</sup>. However, in the most recent cases, published evidence has consisted of *in silico*, early-phase validation. It has been recognized that most recent AI studies are inadequately reported and existing reporting guidelines do not fully cover potential sources of bias specific to AI systems<sup>25</sup>. The welcome emergence of RCTs seeking to evaluate newer interventions based on, or including, an AI component (called ‘AI interventions’ here)<sup>23,26–31</sup> has similarly been met with concerns about the design and reporting<sup>25,32–34</sup>. This has highlighted the need to provide reporting guidance that is ‘fit for purpose’ in this domain.

<sup>1</sup>Moorfields Eye Hospital NHS Foundation Trust, London, UK. <sup>2</sup>Academic Unit of Ophthalmology, Institute of Inflammation and Ageing, University of Birmingham, Birmingham, UK. <sup>3</sup>University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK. <sup>4</sup>Health Data Research UK, London, UK. <sup>5</sup>Birmingham Health Partners Centre for Regulatory Science and Innovation, University of Birmingham, Birmingham, UK. <sup>6</sup>Centre for Patient Reported Outcomes Research, Institute of Applied Health Research, University of Birmingham, Birmingham, UK. <sup>7</sup>Institute of Applied Health Research, University of Birmingham, Birmingham, UK. <sup>8</sup>Centre for Journalology, Clinical Epidemiology Program, Ottawa Hospital Research Institute, Ottawa, Canada. <sup>9</sup>School of Epidemiology and Public Health, Faculty of Medicine, University of Ottawa, Ottawa, Canada. <sup>10</sup>National Institute of Health Research Birmingham Biomedical Research Centre, University of Birmingham and University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK. <sup>11</sup>National Institute of Health Research Applied Research Collaborative West Midlands, Coventry, UK. <sup>12</sup>National Institute of Health Research Surgical Reconstruction and Microbiology Centre, University of Birmingham and University Hospitals Birmingham NHS Foundation Trust, Birmingham, UK. <sup>13</sup>NIHR Biomedical Research Center at Moorfields Eye Hospital NHS Foundation Trust and UCL Institute of Ophthalmology, London, UK. \*A list of authors and their affiliations appears at the end of the paper. ✉e-mail: [a.denniston@bham.ac.uk](mailto:a.denniston@bham.ac.uk)

**Box 1 | Glossary**

**Artificial Intelligence** The science of developing computer systems which can perform tasks normally requiring human intelligence.

**AI intervention** A health intervention that relies upon an AI/ML component to serve its purpose.

**CONSORT** Consolidated Standards of Reporting Trials.

**CONSORT-AI extension item** An additional checklist item to address AI-specific content that is not adequately covered by CONSORT 2010.

**Class-activation map** Class-activation maps are particularly relevant to image classification AI interventions. Class-activation maps are visualizations of the pixels that had the greatest influence on predicted class, by displaying the gradient of the predicted outcome from the model with respect to the input. They are also referred to as ‘saliency maps’ or ‘heat maps’.

**Health outcome** Measured variables in the trial that are used to assess the effects of an intervention.

**Human–AI interaction** The process of how users (humans) interact with the AI intervention, for the AI intervention to function as intended.

**Clinical outcome** Measured variables in the trial which are used to assess the effects of an intervention.

**Delphi study** A research method that derives the collective opinions of a group through a staged consultation of surveys, questionnaires, or interviews, with an aim to reach consensus at the end.

**Development environment** The clinical and operational settings from which the data used for training the model is generated. This includes all aspects of the physical setting (such as geographical location, physical environment), operational setting (such as integration with an electronic record system, installation on a physical device) and clinical setting (such as primary, secondary and/or tertiary care, patient disease spectrum).

**Fine-tuning** Modifications or additional training performed on the AI intervention model, done with the intention of improving its performance.

**Input data** The data that need to be presented to the AI intervention to allow it to serve its purpose.

**Machine learning** A field of computer science concerned with the development of models/algorithms that can solve specific tasks by learning patterns from data, rather than by following explicit rules. It is seen as an approach within the field of AI.

**Operational environment** The environment in which the AI intervention will be deployed, including the infrastructure required to enable the AI intervention to function.

**Output data** The predicted outcome given by the AI intervention based on modeling of the input data. The output data can be presented in different forms, including a classification (including diagnosis, disease severity or stage, or recommendation such as referability), a probability, a class activation map, etc. The output data typically provide additional clinical information and/or trigger a clinical decision.

**Performance error** Instances in which the AI intervention fails to perform as expected. This term can describe different types of failures, and it is up to the investigator to specify what should be considered a performance error, preferably based on prior evidence. This can range from small decreases in accuracy (compared to expected accuracy) to erroneous predictions or the inability to produce an output, in certain cases.

**SPIRIT** Standard Protocol Items: Recommendations for Interventional Trials.

**SPIRIT-AI** An additional checklist item to address AI-specific content that is not adequately covered by SPIRIT 2013.

**SPIRIT-AI elaboration item** Additional considerations to an existing SPIRIT 2013 item when applied to AI interventions.

CONSORT-AI (as part of the SPIRIT-AI and CONSORT-AI initiative) is an international initiative supported by CONSORT and the EQUATOR (Enhancing the Quality and Transparency of Health Research) Network to evaluate the existing CONSORT 2010 statement and to extend or elaborate this guidance where necessary, to support the reporting of clinical trials for AI interventions<sup>35,36</sup>. It is complementary to the SPIRIT-AI statement, which aims to promote high-quality protocol reporting for AI trials. This Consensus Statement describes the methods used to identify and evaluate candidate items and gain consensus. In addition, it also provides the CONSORT-AI checklist, which includes the new extension items and their accompanying explanations.

## Methods

The SPIRIT-AI and CONSORT-AI extensions were simultaneously developed for clinical trial protocols and trial reports. An announcement for the SPIRIT-AI and CONSORT-AI initiative was published in October 2019 (ref. <sup>35</sup>), and the two guidelines were registered as reporting guidelines under development on the EQUATOR library of reporting guidelines in May 2019. Both guidelines were developed in accordance with the EQUATOR Network’s methodological framework<sup>37</sup>. The SPIRIT-AI and CONSORT-AI Steering Group, consisting of 15 international experts, was formed to oversee the

conduct and methodology of the study. Definitions of key terms are provided in the glossary (Box 1).

## Ethical approval

This study was approved by the ethical review committee at the University of Birmingham, UK (ERN\_19-1100). Participant information was provided to Delphi participants electronically before survey completion and before the consensus meeting. Delphi participants provided electronic informed consent, and written consent was obtained from consensus meeting participants.

## Literature review and candidate item generation

An initial list of candidate items for the SPIRIT-AI and CONSORT-AI checklists was generated through review of the published literature and consultation with the Steering Group and known international experts. A search was performed on 13 May 2019 using the terms ‘artificial intelligence’, ‘machine learning’ and ‘deep learning’ to identify existing clinical trials for AI interventions listed within the US National Library of Medicine’s clinical trial registry (ClinicalTrials.gov). There were 316 registered trials, of which 62 were completed and 7 had published results<sup>30,38–43</sup>. Two studies were reported with reference to the CONSORT statement<sup>40,42</sup>, and one study provided an unpublished trial protocol<sup>42</sup>. The Operations

Team (X.L., S.C.R., M.J.C. and A.K.D.) identified AI-specific considerations from these studies and reframed them as candidate reporting items. The candidate items were also informed by findings from a previous systematic review that evaluated the diagnostic accuracy of deep-learning systems for medical imaging<sup>25</sup>. After consultation with the Steering Group and additional international experts ( $n=19$ ), 29 candidate items were generated, 26 of which were relevant for both SPIRIT-AI and CONSORT-AI and 3 of which were relevant only for CONSORT-AI. The Operations Team mapped these items to the corresponding SPIRIT and CONSORT items, revising the wording and providing explanatory text as required to contextualize the items. These items were included in subsequent Delphi surveys.

### Delphi consensus process

In September 2019, 169 key international experts were invited to participate in the online Delphi survey to vote upon the candidate items and suggest additional items. Experts were identified and contacted via the Steering Group and were allowed one round of 'snowball' recruitment in which contacted experts could suggest additional experts. In addition, individuals who made contact following publication of the announcement were included<sup>35</sup>. The Steering Group agreed that individuals with expertise in clinical trials and AI and machine learning (ML), as well as key users of the technology, should be well represented in the consultation. Stakeholders included healthcare professionals, methodologists, statisticians, computer scientists, industry representatives, journal editors, policy makers, health 'informaticists', experts in law and ethics, regulators, patients and funders. Participant characteristics are described in Supplementary Table 1. Two online Delphi surveys were conducted. DelphiManager software (version 4.0), developed and maintained by the COMET (Core Outcome Measures in Effectiveness Trials) initiative, was used to undertake the e-Delphi survey. Participants were given written information about the study and were asked to provide their level of expertise within the fields of (i) AI/ML, and (ii) clinical trials. Each item was presented for consideration (26 for SPIRIT-AI and 29 for CONSORT-AI). Participants were asked to vote on each item using a 9-point scale, as follows: 1–3, not important; 4–6, important but not critical; and 7–9, important and critical. Respondents provided separate ratings for SPIRIT-AI and CONSORT-AI. There was an option to opt out of voting for each item, and each item included space for free text comments. At the end of the Delphi survey, participants had the opportunity to suggest new items. 103 responses were received for the first Delphi round, and 91 responses (88% of participants from round one) were received for the second round. The results of the Delphi survey informed the subsequent international consensus meeting. 12 new items were proposed by the Delphi study participants and were added for discussion at the consensus meeting. Data collected during the Delphi survey were anonymized, and item-level results were presented at the consensus meeting for discussion and voting.

The two-day consensus meeting took place in January 2020 and was hosted by the University of Birmingham, UK, to seek consensus on the content of SPIRIT-AI and CONSORT-AI. 31 international stakeholders from among the Delphi survey participants were invited to discuss the items and vote on their inclusion. Participants were selected to achieve adequate representation from all the stakeholder groups. 41 items were discussed in turn, comprising the 29 items generated in the initial literature review and item-generation phase (26 items relevant to both SPIRIT-AI and CONSORT-AI; 3 items relevant only to CONSORT-AI) and the 12 new items proposed by participants during the Delphi surveys. Each item was presented to the consensus group, alongside its score from the Delphi exercise (median and interquartile ranges) and any comments made by Delphi participants related to that item. Consensus-meeting participants were invited to comment on the importance of each item

and whether the item should be included in the AI extension. In addition, participants were invited to comment on the wording of the explanatory text accompanying each item and the position of each item relative to the SPIRIT 2013 and CONSORT 2010 checklists. After open discussion of each item and the option to adjust wording, an electronic vote took place, with the option to include or exclude the item. An 80% threshold for inclusion was pre-specified and deemed reasonable by the Steering Group to demonstrate majority consensus. Each stakeholder voted anonymously using Turning Point voting pads (Turning Technologies, version 8.7.2.14).

### Checklist pilot

Following the consensus meeting, attendees were given the opportunity to make final comments on the wording and agree that the updated SPIRIT-AI and CONSORT-AI items reflected discussions from the meeting. The Operations Team assigned each item as an extension or elaboration item on the basis of a decision tree and produced a penultimate draft of the SPIRIT-AI and CONSORT-AI checklists (Supplementary Fig. 1). A pilot of the penultimate checklists was conducted with 34 participants to ensure clarity of wording. Experts participating in the pilot included the following: (a) Delphi participants who did not attend the consensus meeting, and (b) external experts who had not taken part in the development process but who had reached out to the Steering Group after the Delphi study commenced. Final changes were made on wording only to improve clarity for readers, by the Operations Team (Supplementary Fig. 2).

### Recommendations

**CONSORT-AI checklist items and explanation.** The CONSORT-AI extension recommends that 14 new checklist items be added to the existing CONSORT 2010 statement (11 extensions and 3 elaborations). These items were considered sufficiently important for clinical-trial reports for AI interventions that they should be routinely reported in addition to the core CONSORT 2010 checklist items. Table 1 lists the CONSORT-AI items.

The 14 items below passed the threshold of 80% for inclusion at the consensus meeting. CONSORT-AI 2a, CONSORT-AI 5 (ii) and CONSORT-AI 19 each resulted from the merging of two items after discussion with the consensus group. CONSORT-AI 4a (i) and (ii) was split into two items for clarity and was voted upon separately. CONSORT-AI 5(iii) did not fulfill the criteria for inclusion on the basis of its initial wording (77% vote to include); however, after extensive discussion and rewording, the consensus group unanimously supported a re-vote, at which point it passed the inclusion threshold (97% to include). The Delphi and voting results for each included and excluded item are described in Supplementary Table 2.

### Title and abstract

**CONSORT-AI 1a,b (i) Elaboration: Indicate that the intervention involves artificial intelligence/machine learning in the title and/or abstract and specify the type of model.** *Explanation.* Indicating in the title and/or abstract of the trial report that the intervention involves a form of AI is encouraged, as it immediately identifies the intervention as an AI/ML intervention and also serves to facilitate indexing and searching of the trial report. The title should be understandable by a wide audience; therefore, a broader umbrella term such as 'artificial intelligence' or 'machine learning' is encouraged. More-precise terms should be used in the abstract, rather than the title, unless they are broadly recognized as being a form of AI/ML. Specific terminology relating to the model type and architecture should be detailed in the abstract.

**CONSORT-AI 1a,b (ii) Elaboration: State the intended use of the AI intervention within the trial in the title and/or abstract.** *Explanation.* Describe the intended use of the AI intervention in the

**Table 1 | CONSORT-AI checklist**

Section		CONSORT 2010 item <sup>a</sup>	CONSORT-AI item	Addressed on page number <sup>b</sup>
Title and abstract				
<b>Title and Abstract</b>	1a	Identification as a randomized trial in the title	CONSORT-AI 1a,b Elaboration	(i) Indicate that the intervention involves artificial intelligence/machine learning in the title and/or abstract and specify the type of model.
	1b	Structured summary of trial design, methods, results, and conclusions (for specific guidance see CONSORT for abstracts)		(ii) State the intended use of the AI intervention within the trial in the title and/or abstract.
Introduction				
<b>Background and objectives</b>	2a	Scientific background and explanation of rationale	CONSORT-AI 2a (i) Extension	Explain the intended use of the AI intervention in the context of the clinical pathway, including its purpose and its intended users (for example, healthcare professionals, patients, public).
	2b	Specific objectives or hypotheses		
Methods				
<b>Trial design</b>	3a	Description of trial design (such as parallel, factorial) including allocation ratio		
	3b	Important changes to methods after trial commencement (such as eligibility criteria), with reasons		
<b>Participants</b>	4a	Eligibility criteria for participants	CONSORT-AI 4a (i) Elaboration	State the inclusion and exclusion criteria at the level of participants.
			CONSORT-AI 4a (ii) Extension	State the inclusion and exclusion criteria at the level of the input data.
	4b	Settings and locations where the data were collected	CONSORT-AI 4b Extension	Describe how the AI intervention was integrated into the trial setting, including any onsite or offsite requirements.
<b>Interventions</b>	5	The interventions for each group with sufficient details to allow replication, including how and when they were actually administered	CONSORT-AI 5 (i) Extension	State which version of the AI algorithm was used.
			CONSORT-AI 5 (ii) Extension	Describe how the input data were acquired and selected for the AI intervention.
			CONSORT-AI 5 (iii) Extension	Describe how poor quality or unavailable input data were assessed and handled.
			CONSORT-AI 5 (iv) Extension	Specify whether there was human-AI interaction in the handling of the input data, and what level of expertise was required of users.
			CONSORT-AI 5 (v) Extension	Specify the output of the AI intervention
			CONSORT-AI 5 (vi) Extension	Explain how the AI intervention's outputs contributed to decision-making or other elements of clinical practice.
<b>Outcomes</b>	6a	Completely defined pre-specified primary and secondary outcome measures, including how and when they were assessed		
	6b	Any changes to trial outcomes after the trial commenced, with reasons		
<b>Sample size</b>	7a	How sample size was determined		
	7b	When applicable, explanation of any interim analyses and stopping guidelines		

Continued



**Table 1 | CONSORT-AI checklist (continued)**

Section	CONSORT 2010 item <sup>a</sup>	CONSORT-AI item	Addressed on page number <sup>b</sup>
Randomization			
<b>Sequence generation</b>	8a	Method used to generate the random allocation sequence	
	8b	Type of randomization; details of any restriction (such as blocking and block size)	
<b>Allocation concealment mechanism</b>	9	Mechanism used to implement the random allocation sequence (such as sequentially numbered containers), describing any steps taken to conceal the sequence until interventions were assigned	
<b>Implementation</b>	10	Who generated the random allocation sequence, who enrolled participants, and who assigned participants to interventions	
<b>Blinding</b>	11a	If done, who was blinded after assignment to interventions (for example, participants, care providers, those assessing outcomes) and how	
	11b	If relevant, description of the similarity of interventions	
<b>Statistical methods</b>	12a	Statistical methods used to compare groups for primary and secondary outcomes	
	12b	Methods for additional analyses, such as subgroup analyses and adjusted analyses	
Results			
<b>Participant flow (a diagram is strongly recommended)</b>	13a	For each group, the numbers of participants who were randomly assigned, received intended treatment, and were analyzed for the primary outcome	
	13b	For each group, losses and exclusions after randomization, together with reasons	
<b>Recruitment</b>	14a	Dates defining the periods of recruitment and follow-up	
	14b	Why the trial ended or was stopped	
<b>Baseline data</b>	15	A table showing baseline demographic and clinical characteristics for each group	
<b>Numbers analyzed</b>	16	For each group, number of participants (denominator) included in each analysis and whether the analysis was by original assigned groups	
<b>Outcomes and estimation</b>	17a	For each primary and secondary outcome, results for each group, and the estimated effect size and its precision (such as 95% confidence interval)	
	17b	For binary outcomes, presentation of both absolute and relative effect sizes is recommended	
<b>Ancillary analyses</b>	18	Results of any other analyses performed, including subgroup analyses and adjusted analyses, distinguishing pre-specified from exploratory	

Continued

**Table 1 | CONSORT-AI checklist (continued)**

Section	CONSORT 2010 item <sup>a</sup>	CONSORT-AI item	Addressed on page number <sup>b</sup>
<b>Harms</b>	19	All important harms or unintended effects in each group (for specific guidance see CONSORT for harms)	CONSORT-AI 19 Extension Describe results of any analysis of performance errors and how errors were identified, where applicable. If no such analysis was planned or done, justify why not.
Discussion			
<b>Limitations</b>	20	Trial limitations, addressing sources of potential bias, imprecision, and, if relevant, multiplicity of analyses	
<b>Generalizability</b>	21	Generalizability (external validity, applicability) of the trial findings	
<b>Interpretation</b>	22	Interpretation consistent with results, balancing benefits and harms, and considering other relevant evidence	
Other Information			
<b>Registration</b>	23	Registration number and name of trial registry	
<b>Protocol</b>	24	Where the full trial protocol can be accessed, if available	
<b>Funding</b>	25	Sources of funding and other support (such as supply of drugs), role of funders	CONSORT-AI 25 Extension State whether and how the AI intervention and/or its code can be accessed, including any restrictions to access or re-use.

<sup>a</sup>We strongly recommend reading this statement in conjunction with the CONSORT 2010 Explanation and Elaboration for important clarifications on all the items. <sup>b</sup>Indicates page numbers to be completed by authors during protocol development.

trial report title and/or abstract. This should describe the purpose of the AI intervention and the disease context<sup>26,44</sup>. Some AI interventions may have multiple intended uses, or the intended use may evolve over time. Therefore, documenting this allows readers to understand the intended use of the algorithm at the time of the trial.

## Introduction

**CONSORT-AI 2a (i) Extension: Explain the intended use for the AI intervention in the context of the clinical pathway, including its purpose and its intended users (for example, healthcare professionals, patients, public).** *Explanation.* In order to clarify how the AI intervention is intended to fit into a clinical pathway, a detailed description of its role should be included in the background of the trial report. AI interventions may be designed to interact with different users, including healthcare professionals, patients and the public, and their roles can be wide-ranging (for example, the same AI intervention could theoretically be replacing, augmenting or adjudicating components of clinical decision-making). Clarifying the intended use of the AI intervention and its intended user helps readers understand the purpose for which the AI intervention was evaluated in the trial.

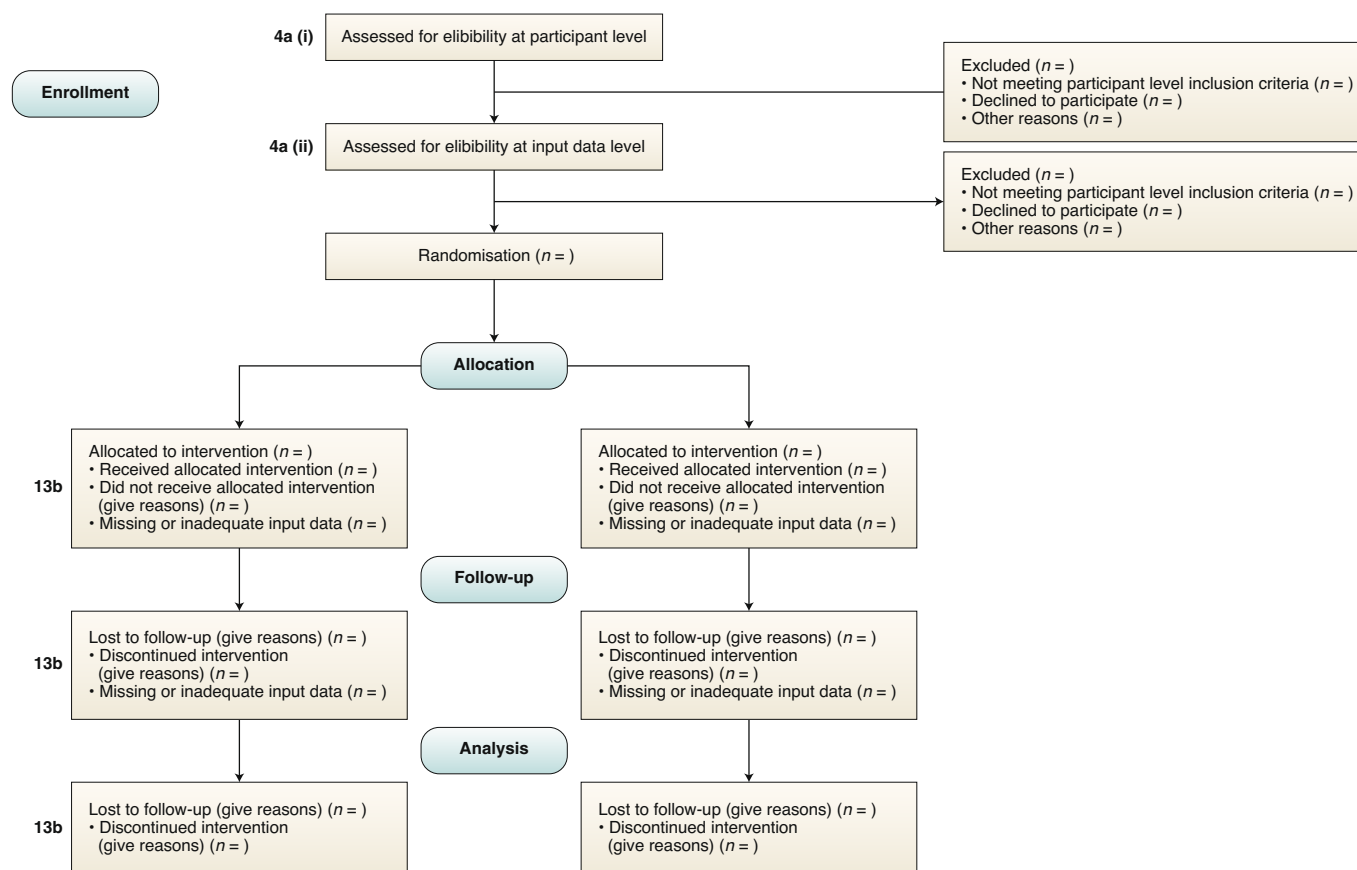
## Methods

**CONSORT-AI 4a (i) Elaboration: State the inclusion and exclusion criteria at the level of participants.** *Explanation.* The inclusion and exclusion criteria should be defined at the participant level as per usual practice in non-AI interventional trial reports (Fig. 1). This is distinct from the inclusion and exclusion criteria made at the input-data level, which is addressed in item 4a (ii).

**CONSORT-AI 4a (ii) Extension: State the inclusion and exclusion criteria at the level of the input data.** *Explanation.* ‘Input data’ refers to the data required by the AI intervention to serve its

purpose (for example, for a breast-cancer diagnostic system, the input data could be the unprocessed or vendor-specific post-processing mammography scan upon which a diagnosis is being made; for an early-warning system, the input data could be physiological measurements or laboratory results from the electronic health record). The trial report should pre-specify if there were minimum requirements for the input data (such as image resolution, quality metrics or data format) that determined pre-randomization eligibility. It should specify when, how and by whom this was assessed. For example, if a participant met the eligibility criteria for lying flat for a CT scan as per item 4a (i), but the scan quality was compromised (for any given reason) to such a level that it was deemed unfit for use by the AI system, this should be reported as an exclusion criterion at the input-data level. Note that where input data are acquired after randomization, any exclusion is considered to be from the analysis, not from enrollment (CONSORT item 13b and Fig. 1).

**CONSORT-AI 4b Extension: Describe how the AI intervention was integrated into the trial setting, including any onsite or offsite requirements.** *Explanation.* There are limitations to the generalizability of AI algorithms, one of which is when they are used outside of their development environment<sup>45,46</sup>. AI systems are dependent on their operational environment, and the report should provide details of the hardware and software requirements to allow technical integration of the AI intervention at each study site. For example, it should be stated if the AI intervention required vendor-specific devices, if there was specialized computing hardware at each site, or if the site had to support cloud integration, particularly if this was vendor specific. If any changes to the algorithm were required at each study site as part of the implementation procedure (such as fine-tuning the algorithm on local data), then this process should also be clearly described.



**Fig. 1 | CONSORT 2010 flow diagram — adapted for AI clinical trials.** CONSORT-AI 4a (i): State the inclusion and exclusion criteria at the level of participants. CONSORT-AI 4a (ii): State the inclusion and exclusion criteria at the level of the input data. CONSORT 13b (core CONSORT item): For each group, losses and exclusions after randomization, together with reasons.

**CONSORT-AI 5 (i) Extension: State which version of the AI algorithm was used.** *Explanation.* Similar to other forms of software as a medical device, AI systems are likely to undergo multiple iterations and updates during their lifespan. It is therefore important to specify which version of the AI system was used in the clinical trial, whether this is the same as the version evaluated in previous studies that have been used to justify the study rationale, and whether the version changed during the conduct of the trial. If applicable, the report should describe what has changed between the relevant versions and the rationales for the changes. Where available, the report should include a regulatory marking reference, such as a unique device identifier, that requires a new identifier for updated versions of the device<sup>47</sup>.

**CONSORT-AI 5 (ii) Extension: Describe how the input data were acquired and selected for the AI intervention.** *Explanation.* The measured performance of any AI system may be critically dependent on the nature and quality of the input data<sup>48</sup>. A description of the input-data handling, including acquisition, selection and pre-processing before analysis by the AI system, should be provided. Completeness and transparency of this description is integral to the replicability of the intervention beyond the clinical trial in real-world settings. It also helps readers identify whether input-data-handling procedures were standardized across trial sites.

**CONSORT-AI 5 (iii) Extension: Describe how poor-quality or unavailable input data were assessed and handled.** *Explanation.* As with CONSORT-AI 4a (ii), ‘input data’ refers to the data required by the AI intervention to serve its purpose. As discussed in item

4a (ii), the performance of AI systems may be compromised as a result of poor quality or missing input data<sup>49</sup> (for example, excessive movement artifact on an electrocardiogram). The trial report should report the amount of missing data, as well as how this was identified and handled. The report should also specify if there was a minimum standard required for the input data and, where this standard was not achieved, how this was handled (including the impact on, or any changes to, the participant care pathway).

Poor quality or unavailable data can also affect non-AI interventions. For example, sub-optimal quality of a scan could affect a radiologist’s ability to interpret it and make a diagnosis. It is therefore important that this information is reported equally in the control intervention, where relevant. If this minimum quality standard was different from the inclusion criteria for input data used to assess eligibility pre-randomization, this should be stated.

**CONSORT-AI 5 (iv) Extension: Specify whether there was human–AI interaction in the handling of the input data, and what level of expertise was required of users.** *Explanation.* A description of the human–AI interface and the requirements for successful interaction when input data are handled should be provided — for example, clinician-led selection of regions of interest from a histology slide that is then interpreted by an AI diagnostic system<sup>50</sup>, or an endoscopist’s selection of a colonoscopy video clips as input data for an algorithm designed to detect polyps<sup>28</sup>. A description of any user training provided and instructions for how users should handle the input data provides transparency and replicability of trial procedures. Poor clarity on the human–AI interface may lead to lack of a standard approach and may carry ethical



implications, particularly in the event of harm<sup>51,52</sup>. For example, it may become unclear whether an error case occurred due to human deviation from the instructed procedure, or if it was an error made by the AI system.

**CONSORT-AI 5 (v) Extension: Specify the output of the AI intervention.** *Explanation.* The output of the AI intervention should be clearly specified in the trial report. For example, an AI system may output a diagnostic classification or probability, a recommended action, an alarm alerting to an event, an instigated action in a closed-loop system (such as titration of drug infusions) or another output. The nature of the AI intervention's output has direct implications on its usability and how it may lead to downstream actions and outcomes.

**CONSORT-AI 5 (vi) Extension: Explain how the AI intervention's outputs contributed to decision-making or other elements of clinical practice.** *Explanation.* Since health outcomes may also critically depend on how humans interact with the AI intervention, the report should explain how the outputs of the AI system were used to contribute to decision-making or other elements of clinical practice. This should include adequate description of downstream interventions that can affect outcomes. As with CONSORT-AI 5 (iv), any effects of human–AI interaction on the outputs should be described in detail, including the level of expertise required to understand the outputs and any training and/or instructions provided for this purpose. For example, a skin cancer detection system that produced a percentage likelihood as its output should be accompanied by an explanation of how this output was interpreted and acted upon by the user, specifying both the intended pathways (for example, skin lesion excision if the diagnosis is positive) and the thresholds for entry to these pathways (for example, skin lesion excision if the diagnosis is positive and the probability is >80%). The information produced by comparator interventions should be similarly described, alongside an explanation of how such information was used to arrive at clinical decisions on patient management, where relevant. Any discrepancy in how decision-making occurred versus how it was intended to occur (that is, as specified in the trial protocol) should be reported.

## Results

**CONSORT-AI 19 Extension: Describe results of any analysis of performance errors and how errors were identified, where applicable. If no such analysis was planned or done, explain why not.** *Explanation.* Reporting performance errors and failure case analysis is especially important for AI interventions. AI systems can make errors that may be hard to foresee but that, if allowed to be deployed at scale, could have catastrophic consequences<sup>53</sup>. Therefore, reporting cases of error and defining risk-mitigation strategies are important for informing when, and for which populations, the intervention can be safely implemented. The results of any performance-error analysis should be reported and the implications of the results should be discussed.

## Other information

**CONSORT-AI 25 Extension: State whether and how the AI intervention and/or its code can be accessed, including any restrictions to access or re-use.** *Explanation.* The trial report should make it clear whether and how the AI intervention and/or its code can be accessed or re-used. This should include details about the license and any restrictions to access.

## Discussion

CONSORT-AI is a new reporting-guideline extension developed through international multi-stakeholder consensus. It aims to promote transparent reporting of AI intervention trials and

is intended to facilitate critical appraisal and evidence synthesis. The extension items added in CONSORT-AI address a number of issues specific to the implementation and evaluation of AI interventions, which should be considered alongside the core CONSORT 2010 checklist and other CONSORT extensions<sup>54</sup>. It is important to note that these are minimum requirements and there may be value in including additional items not included in the checklists in the report or in supplementary materials (Supplementary Table 2).

In both CONSORT-AI and its companion project SPIRIT-AI, a major emphasis was the addition of several new items related to the intervention itself and its application in the clinical context. Items 5 (i)–5 (vi) were added to address AI-specific considerations in descriptions of the intervention. Specific recommendations were made pertinent to AI systems related to algorithm version, input and output data, integration into trial settings, expertise of the users and protocol for acting upon the AI system's recommendations. It was agreed that these details are critical for independent evaluation or replication of the trial. Journal editors reported that despite the importance of these items, they are currently often missing from trial reports at the time of submission for publication, which provides further weight for their inclusion as specifically listed extension items.

A recurrent focus of the Delphi comments and consensus group discussion was the safety of AI systems. This was in recognition that AI systems, unlike other health interventions, can unpredictably yield errors that are not easily detectable or explainable by human judgement. For example, changes to medical imaging that are invisible, or appear random, to the human eye may change the likelihood of the diagnostic output entirely<sup>55,56</sup>. The concern is that given the theoretical ease with which AI systems could be deployed at scale, any unintended harmful consequences could be catastrophic. CONSORT-AI item 19, which requires specification of any plans to analyze performance errors, was added to emphasize the importance of anticipating systematic errors made by the algorithm and their consequences. Beyond this, investigators should also be encouraged to explore differences in performance and error rates across population subgroups. It has been shown that AI systems may be systematically biased toward different outputs, which may lead to different or even unfair treatment, on the basis of extant features<sup>53,57–59</sup>.

The topic of 'continuously evolving' AI systems (also known as 'continuously adapting' or 'continuously learning' AI systems) was discussed at length during the consensus meeting, but it was agreed that this be excluded from CONSORT-AI. These are AI systems with the ability to continuously train on new data, which may cause changes in performance over time. The group noted that, while of interest, this field is relatively early in its development without tangible examples in healthcare applications, and that it would not be appropriate for it to be included in CONSORT-AI at this stage<sup>60</sup>. This topic will be monitored and will be revisited in future iterations of CONSORT-AI. It is worth noting that incremental software changes, whether continuous or iterative, intentional or unintentional, could have serious consequences on safety performance after deployment. It is therefore of vital importance that such changes be documented and identified by software version and that a robust post-deployment surveillance plan is in place.

This study is set in the current context of AI in health; therefore, several limitations should be noted. First, there are relatively few published interventional trials in the field of AI for healthcare; therefore, the discussions and decisions made during this study were not always supported by existing examples of completed trials. This arises from our stated aim of addressing the issues of poor reporting in this field as early as possible, recognizing the strong drivers in the field and the specific challenges of study design

and reporting for AI. As the science and study of AI evolves, we welcome collaboration with investigators to co-evolve these reporting standards to ensure their continued relevance. Second, the literature search of AI RCTs used terminology such as ‘artificial intelligence’, ‘machine learning’ and ‘deep learning’, but not terms such as ‘clinical decision support systems’ or ‘expert systems’, which were more commonly used in the 1990s for technologies underpinned by AI systems and share risks similar to those of recent examples<sup>61</sup>. It is likely that such systems, if published today, would be indexed under ‘artificial intelligence’ or ‘machine learning’; however, clinical decision support systems were not actively discussed during this consensus process. Third, the initial candidate items list was generated by a relatively small group of experts consisting of SteeringGroup members and additional international experts; however, additional items from the wider Delphi group were taken forward for consideration by the consensus group, and no new items were suggested during the consensus meeting or post-meeting evaluation.

As with the CONSORT statement, the CONSORT-AI extension is intended as a minimum reporting guidance, and there are additional AI-specific considerations for trial reports that may warrant consideration (Supplementary Table 2). This extension is aimed particularly at investigators and readers reporting or appraising clinical trials; however, it may also serve as useful guidance for developers of AI interventions in earlier validation stages of an AI system. Investigators seeking to report studies developing and validating the diagnostic and predictive properties of AI models should refer to TRIPOD-ML (Transparent Reporting of a Multivariable Prediction Model for Individual Prognosis or Diagnosis–Machine Learning) and STARD-AI (Standards for Reporting Diagnostic Accuracy Studies–Artificial Intelligence), both of which are currently under development<sup>32,62</sup>. Other potentially relevant guidelines, which are agnostic to study design, are registered with the EQUATOR Network<sup>63</sup>. The CONSORT-AI extension is expected to encourage careful early planning of AI interventions for clinical trials and this, in conjunction with SPIRIT-AI, should help to improve the quality of trials for AI interventions. The development of the CONSORT-AI guidance does not include additional items within the discussion section of trial reports. The guidance provided by CONSORT 2010 on trial limitations, generalizability and interpretation were deemed to be translatable to trials for AI interventions.

There is also recognition that AI is a rapidly evolving field, and there will be the need to update CONSORT-AI as the technology, and newer applications for it, develop. Currently, most applications of AI involve disease detection, diagnosis and triage, and this is likely to have influenced the nature and prioritization of items within CONSORT-AI. As wider applications that utilize ‘AI as therapy’ emerge, it will be important to continue to evaluate CONSORT-AI in the light of such studies. Additionally, advances in computational techniques and the ability to integrate them into clinical workflows will bring new opportunities for innovation that benefits patients. However, they may be accompanied by new challenges around study design and reporting. In order to ensure transparency, minimize potential biases and promote the trustworthiness of the results and the extent to which they may be generalizable, the SPIRIT-AI and CONSORT-AI Steering Group will continue to monitor the need for updates.

### Data availability

Data requests should be made to the corresponding author and release will be subject to consideration by the SPIRIT-AI and CONSORT-AI Steering Group.

Received: 24 April 2020; Accepted: 23 July 2020;  
Published online: 9 September 2020

### References

- Sibbald, B. & Roland, M. Understanding controlled trials. Why are randomised controlled trials important? *Br. Med. J.* **316**, 201 (1998).
- Peto, R., Collins, R. & Gray, R. Large-scale randomized evidence: large, simple trials and overviews of trials. *J. Clin. Epidemiol.* **48**, 23–40 (1995).
- Jüni, P., Altman, D. G. & Egger, M. Systematic reviews in health care: assessing the quality of controlled clinical trials. *Br. Med. J.* **323**, 42–46 (2001).
- Schulz, K. F., Chalmers, I., Hayes, R. J. & Altman, D. G. Empirical evidence of bias. Dimensions of methodological quality associated with estimates of treatment effects in controlled trials. *J. Am. Med. Assoc.* **273**, 408–412 (1995).
- Moher, D. et al. CONSORT 2010 explanation and elaboration: updated guidelines for reporting parallel group randomised trials. *Br. Med. J.* **340**, c869 (2010).
- Moher, D., Jones, A. & Lepage, L. CONSORT Group (Consolidated Standards for Reporting of Trials) Use of the CONSORT statement and quality of reports of randomized trials: a comparative before-and-after evaluation. *J. Am. Med. Assoc.* **285**, 1992–1995 (2001).
- Glasziou, P. et al. Reducing waste from incomplete or unusable reports of biomedical research. *Lancet* **383**, 267–276 (2014).
- Boutron, I., Altman, D. G., Moher, D., Schulz, K. F. & Ravaud, P. CONSORT statement for randomized trials of nonpharmacologic treatments: a 2017 update and a CONSORT extension for nonpharmacologic trial abstracts. *Ann. Intern. Med.* **167**, 40–47 (2017).
- Hopewell, S. et al. CONSORT for reporting randomised trials in journal and conference abstracts. *Lancet* **371**, 281–283 (2008).
- MacPherson, H. et al. Revised STAndards for Reporting Interventions in Clinical Trials of Acupuncture (STRICTA): extending the CONSORT statement. *PLoS Med.* **7**, e1000261 (2010).
- Gagnier, J. J. et al. Reporting randomized, controlled trials of herbal interventions: an elaborated CONSORT statement. *Ann. Intern. Med.* **144**, 364–367 (2006).
- Cheng, C.-W. et al. CONSORT extension for Chinese herbal medicine formulas 2017: recommendations, explanation, and elaboration. *Ann. Intern. Med.* **167**, 112–121 (2017).
- Calvert, M. et al. Reporting of patient-reported outcomes in randomized trials: the CONSORT PRO extension. *J. Am. Med. Assoc.* **309**, 814–822 (2013).
- He, J. et al. The practical implementation of artificial intelligence technologies in medicine. *Nat. Med.* **25**, 30–36 (2019).
- McKinney, S. M. et al. International evaluation of an AI system for breast cancer screening. *Nature* **577**, 89–94 (2020).
- Abramoff, M. D. et al. Improved automated detection of diabetic retinopathy on a publicly available dataset through integration of deep learning. *Invest. Ophthalmol. Vis. Sci.* **57**, 5200–5206 (2016).
- De Fauw, J. et al. Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nat. Med.* **24**, 1342–1350 (2018).
- Esteva, A. et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nature* **542**, 115–118 (2017).
- Rajpurkar, P. et al. Deep learning for chest radiograph diagnosis: A retrospective comparison of the CheXNeXt algorithm to practicing radiologists. *PLoS Med.* **15**, e1002686 (2018).
- Fleuren, L. M. et al. Machine learning for the prediction of sepsis: a systematic review and meta-analysis of diagnostic test accuracy. *Intensive Care Med.* **46**, 383–400 (2020).
- Yim, J. et al. Predicting conversion to wet age-related macular degeneration using deep learning. *Nat. Med.* **26**, 892–899 (2020).
- Kim, H., Goo, J. M., Lee, K. H., Kim, Y. T. & Park, C. M. Preoperative CT-based deep learning model for predicting disease-free survival in patients with lung adenocarcinomas. *Radiology* **296**, 216–224 (2020).
- Wang, P. et al. Real-time automatic detection system increases colonoscopic polyp and adenoma detection rates: a prospective randomised controlled study. *Gut* **68**, 1813–1819 (2019).
- Tyler, N. S. et al. An artificial intelligence decision support system for the management of type 1 diabetes. *Nat. Metab.* **2**, 612–619 (2020).
- Liu, X. et al. A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: a systematic review and meta-analysis. *Lancet Digit. Health* **1**, e271–e297 (2019).
- Wijnberge, M. et al. Effect of a machine learning-derived early warning system for intraoperative hypotension vs standard care on depth and duration of intraoperative hypotension during elective noncardiac surgery: the HYPE randomized clinical trial. *J. Am. Med. Assoc.* **323**, 1052–1060 (2020).
- Gong, D. et al. Detection of colorectal adenomas with a real-time computer-aided system (ENDOANGEL): a randomised controlled study. *Lancet Gastroenterol. Hepatol.* **5**, 352–361 (2020).
- Wang, P. et al. Effect of a deep-learning computer-aided detection system on adenoma detection during colonoscopy (CADE-DB trial): a double-blind randomised study. *Lancet Gastroenterol. Hepatol.* **5**, 343–351 (2020).

29. Wu, L. et al. Randomised controlled trial of WISENSE, a real-time quality improving system for monitoring blind spots during esophagogastroduodenoscopy. *Gut* **68**, 2161–2169 (2019).
30. Lin, H. et al. Diagnostic efficacy and therapeutic decision-making capacity of an artificial intelligence platform for childhood cataracts in eye clinics: a multicentre randomized controlled trial. *EClinicalMedicine* **9**, 52–59 (2019).
31. Su, J.-R. et al. Impact of a real-time automatic quality control system on colorectal polyp and adenoma detection: a prospective randomized controlled study (with videos). *Gastrointest. Endosc.* **91**, 415–424.e4 (2020).
32. Collins, G. S. & Moons, K. G. M. Reporting of artificial intelligence prediction models. *Lancet* **393**, 1577–1579 (2019).
33. Gregory, J., Welliver, S. & Chong, J. Top 10 reviewer critiques of radiology artificial intelligence (AI) articles: qualitative thematic analysis of reviewer critiques of machine learning/deep learning manuscripts submitted to JMIR. *J. Magn. Reson. Imaging* **52**, 248–254 (2020).
34. Nagendran, M. et al. Artificial intelligence versus clinicians: systematic review of design, reporting standards, and claims of deep learning studies. *Br. Med. J.* **368**, m689 (2020).
35. CONSORT-AI and SPIRIT-AI Steering Group. Reporting guidelines for clinical trials evaluating artificial intelligence interventions are needed. *Nat. Med.* **25**, 1467–1468 (2019).
36. Liu, X., Faes, L., Calvert, M. J. & Denniston, A. K. Extension of the CONSORT and SPIRIT statements. *Lancet* **394**, 1225 (2019).
37. Moher, D., Schulz, K. F., Simera, I. & Altman, D. G. Guidance for developers of health research reporting guidelines. *PLoS Med.* **7**, e1000217 (2010).
38. Caballero-Ruiz, E. et al. A web-based clinical decision support system for gestational diabetes: Automatic diet prescription and detection of insulin needs. *Int. J. Med. Inform.* **102**, 35–49 (2017).
39. Kim, T. W. B., Gay, N., Khemka, A. & Garino, J. Internet-based exercise therapy using algorithms for conservative treatment of anterior knee pain: a pragmatic randomized controlled trial. *JMIR Rehabil. Assist. Technol.* **3**, e12 (2016).
40. Labovitz, D. L., Shafner, L., Reyes Gil, M., Virmani, D. & Hanina, A. Using artificial intelligence to reduce the risk of nonadherence in patients on anticoagulation therapy. *Stroke* **48**, 1416–1419 (2017).
41. Nicolae, A. et al. Evaluation of a machine-learning algorithm for treatment planning in prostate low-dose-rate brachytherapy. *Int. J. Radiat. Oncol. Biol. Phys.* **97**, 822–829 (2017).
42. Voss, C. et al. Effect of wearable digital intervention for improving socialization in children with autism spectrum disorder: a randomized clinical trial. *JAMA Pediatr.* **173**, 446–454 (2019).
43. Mendes-Soares, H. et al. Assessment of a personalized approach to predicting postprandial glycemic responses to food among individuals without diabetes. *JAMA Netw. Open* **2**, e188102 (2019).
44. Choi, K. J. et al. Development and validation of a deep learning system for staging liver fibrosis by using contrast agent-enhanced CT images in the liver. *Radiology* **289**, 688–697 (2018).
45. Kelly, C. J., Karthikesalingam, A., Suleyman, M., Corrado, G. & King, D. Key challenges for delivering clinical impact with artificial intelligence. *BMC Med.* **17**, 195 (2019).
46. Pooch, E. H. P., Ballester, P. L. & Barros, R. C. Can we trust deep learning models diagnosis? The impact of domain shift in chest radiograph classification. *arXiv* <http://arxiv.org/abs/1909.01940> (2019).
47. International Medical Device Regulators Forum. Unique Device Identification System (UDI System) Application Guide. <http://www.imdrf.org/documents/documents.asp> (2019).
48. Sabottke, C. F. & Spieler, B. M. The effect of image resolution on deep learning in radiography. *Radiology: Artif. Intell.* **2**, e190015 (2020).
49. Heaven, D. Why deep-learning AIs are so easy to fool. *Nature* **574**, 163–166 (2019).
50. Kiani, A. et al. Impact of a deep learning assistant on the histopathologic classification of liver cancer. *NPJ Digit. Med.* **3**, 23 (2020).
51. Wiens, J. et al. Do no harm: a roadmap for responsible machine learning for health care. *Nat. Med.* **25**, 1337–1340 (2019).
52. Habli, I., Lawton, T. & Porter, Z. Artificial intelligence in health care: accountability and safety. *Bull. World Health Organ.* [https://www.who.int/bulletin/online\\_first/BLT.19.237487.pdf](https://www.who.int/bulletin/online_first/BLT.19.237487.pdf) (2020).
53. Oakden-Rayner, L., Dunnmon, J., Carneiro, G. & Ré, C. Hidden stratification causes clinically meaningful failures in machine learning for medical imaging. *arXiv* <http://arxiv.org/abs/1909.12475> (2019).
54. CONSORT. Extensions of the CONSORT Statement. <http://www.consort-statement.org/extensions> (accessed 24 March 2020).
55. Zech, J. R. et al. Confounding variables can degrade generalization performance of radiological deep learning models. *arXiv* <http://arxiv.org/abs/1807.00431> (2018).
56. Finlayson, S. G. et al. Adversarial attacks on medical machine learning. *Science* **363**, 1287–1289 (2019).
57. Adamson, A. S. & Smith, A. Machine learning and health care disparities in dermatology. *JAMA Dermatol.* **154**, 1247–1248 (2018).
58. Zou, J. & Schiebinger, L. AI can be sexist and racist - it's time to make it fair. *Nature* **559**, 324–326 (2018).
59. Chen, I. Y., Joshi, S. & Ghassemi, M. Treating health disparities with artificial intelligence. *Nat. Med.* **26**, 16–17 (2020).
60. Lee, C. S. & Lee, A. Y. Clinical applications of continual learning machine learning. *Lancet Digit. Health* **2**, e279–e281 (2020).
61. Sutton, R. T. et al. An overview of clinical decision support systems: benefits, risks, and strategies for success. *NPJ Digit. Med.* **3**, 17 (2020).
62. Sounderajah, V. et al. Developing specific reporting guidelines for diagnostic accuracy studies assessing AI interventions: The STARD-AI Steering Group. *Nat. Med.* **26**, 807–808 (2020).
63. Talmon, J. et al. STARE-HI-Statement on reporting of evaluation studies in Health Informatics. *Int. J. Med. Inform.* **78**, 1–9 (2009).

## Acknowledgements

We thank the participants who were involved in the Delphi study and Pilot study (Supplementary Note); E. Marston (University of Birmingham, UK) for strategic support; and C. Radovanovic (University Hospitals Birmingham NHS Foundation Trust, UK) and A. Walker (University of Birmingham, UK) for administrative support. The views expressed in this publication are those of the authors, Delphi participants and stakeholder participants and may not represent the views of the broader stakeholder group or host institution. This work was funded by a Wellcome Trust Institutional Strategic Support Fund: Digital Health Pilot Grant, Research England (part of UK Research and Innovation), Health Data Research UK and the Alan Turing Institute. The study was sponsored by the University of Birmingham, UK. The study funders and sponsors had no role in the design and conduct of the study; collection, management, analysis and interpretation of the data; preparation, review or approval of the manuscript; or decision to submit the manuscript for publication. M.J.C. is a National Institute for Health Research (NIHR) Senior Investigator and receives funding from the National Institute for Health Research (NIHR) Birmingham Biomedical Research Centre; the NIHR Surgical Reconstruction and Microbiology Research Centre and NIHR ARC West Midlands at the University of Birmingham and University Hospitals Birmingham NHS Foundation Trust; Health Data Research UK; Innovate UK (part of UK Research and Innovation); the Health Foundation; Macmillan Cancer Support; and UCB Pharma. A.D. and J.D. are also NIHR Senior Investigators. The views expressed in this article are those of the author(s) and not necessarily those of the NIHR, or the Department of Health and Social Care. S.J.V. receives funding from the Engineering and Physical Sciences Research Council, UK Research and Innovation (UKRI), Accenture, Warwick Impact Fund, Health Data Research UK and European Regional Development Fund. S.R. is an employee of the Medical Research Council (UKRI). D.M. is supported by a University of Ottawa Research Chair. M.K.E. is supported by the U.S. Food and Drug Administration (FDA), and D.P. is supported in part by the Office of the Director at the National Library of Medicine (NLM), US National Institutes of Health (NIH). A.B. is supported by a National Institutes of Health (NIH) award 7K01HL141771-02. This article may not be consistent with NIH and/or FDA's views or policies. It reflects only the views and opinions of the authors.

## Author contributions

Concept and design, and acquisition, analysis and interpretation of data, all authors; drafting of the manuscript, X.L., S.C.R., D.M., M.J.C. and A.K.D.; obtaining of funding, M.J.C., C.Y., C.H. and A.K.D. The SPIRIT-AI and CONSORT-AI Working Group consists of two groups that have been key in the development of the guidelines: the SPIRIT-AI and CONSORT-AI Steering Group, which was responsible for overseeing the consensus process and guidelines development methodology (Alastair K. Denniston, An-Wen Chan, Ara Darzi, Christopher Holmes, Christopher Yau, David Moher, Hután Ashrafian, Jonathan J. Deeks, Lavinia Ferrante di Ruffano, Livia Faes, Melanie J. Calvert, Pearse A. Keane, Samantha Cruz Rivera, Sebastian J. Vollmer and Xiaoxuan Liu); and the SPIRIT-AI and CONSORT-AI Consensus Group, which was responsible for reaching consensus on the content and wording of the items within the checklists (Aaron Y. Lee, Adrian Jonas, Andre Esteve, Andrew L. Beam, An-Wen Chan, Maria Beatrice Panico, Cecilia S. Lee, Charlotte Haug, Christopher J. Kelly, Christopher Yau, Cynthia Mulrow, Cyrus Espinoza, David Moher, Dina Paltoo, Elaine Manna, Gary Price, Gary S. Collins, Hugh Harvey, James Matcham, Joao Monteiro, John Fletcher, M. Khair ElZarrad, Lavinia Ferrante Di Ruffano, Luke Oakden-Rayner, Melanie J. Calvert, Melissa McCradden, Pearse A. Keane, Richard Savage, Robert Golub, Rupa Sarkar and Samuel Rowley).

## Competing interests

M.J.C. has received personal fees from Astellas, Takeda, Merck, Daiichi Sankyo, Glaukos, GlaxoSmithKline and the Patient-Centered Outcomes Research Institute (PCORI) outside the submitted work. P.A.K. is a consultant for DeepMind Technologies, Roche, Novartis and Apellis, and has received speaker fees or travel support from Bayer, Allergan, Topcon and Heidelberg Engineering. C.J.K. is an employee of Google and owns Alphabet stock. A.E. is an employee of Salesforce CRM. R.S. is an employee of Pinpoint Science. J. Matcham was an employee of AstraZeneca at the time of this study. J. Monteiro is Chief Editor of the journal *Nature Medicine*; he has recused himself from any aspect of decision-making on this manuscript and played no part in the assignment of this manuscript to in-house editors or peer reviewers, and was also separated and blinded from the editorial process from submission inception to decision.



# Additional information

**Supplementary information** Supplementary information is available for this paper at <https://doi.org/10.1038/s41591-020-1034-x>.

**Correspondence and requests for materials** should be addressed to A.K.D.

**Peer review information** Javier Carmona was the primary editor on this article, and managed its editorial process and peer review in collaboration with the rest of the editorial team.

**Reprints and permissions information** is available at [www.nature.com/reprints](http://www.nature.com/reprints).



**Open Access** This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020

## SPIRIT-AI and CONSORT-AI Steering Group

**Alastair K. Denniston<sup>2,3,4,5,6,13</sup>, An-Wen Chan<sup>14</sup>, Ara Darzi<sup>15,16</sup>, Christopher Holmes<sup>17,18</sup>, Christopher Yau<sup>17,19</sup>, David Moher<sup>8,9</sup>, Hutan Ashrafian<sup>15,16</sup>, Jonathan J. Deeks<sup>7,10</sup>, Lavinia Ferrante di Ruffano<sup>7</sup>, Livia Faes<sup>20</sup>, Melanie J. Calvert<sup>4,5,6,7,10,11,12</sup>, Pearse A. Keane<sup>13</sup>, Samantha Cruz Rivera<sup>5,6,7</sup>, Sebastian J. Vollmer<sup>17,21</sup> and Xiaoxuan Liu<sup>1,2,3,4,5</sup>**

## SPIRIT-AI and CONSORT-AI Consensus Group

**Aaron Y. Lee<sup>22</sup>, Adrian Jonas<sup>23</sup>, Andre Esteva<sup>24</sup>, Andrew L. Beam<sup>25</sup>, An-Wen Chan<sup>14</sup>, Maria Beatrice Panico<sup>26</sup>, Cecilia S. Lee<sup>22</sup>, Charlotte Haug<sup>27</sup>, Christopher J. Kelly<sup>28</sup>, Christopher Yau<sup>17,19</sup>, Cynthia Mulrow<sup>29</sup>, Cyrus Espinoza<sup>30</sup>, John Fletcher<sup>31</sup>, David Moher<sup>8,9</sup>, Dina Paltoo<sup>32</sup>, Elaine Manna<sup>33</sup>, Gary Price<sup>34</sup>, Gary S. Collins<sup>35</sup>, Hugh Harvey<sup>36</sup>, James Matcham<sup>37</sup>, Joao Monteiro<sup>38</sup>, M. Khair ElZarrad<sup>39</sup>, Lavinia Ferrante di Ruffano<sup>7</sup>, Luke Oakden-Rayner<sup>40</sup>, Melanie J. Calvert<sup>4,5,6,7,10,11,12</sup>, Melissa McCradden<sup>41</sup>, Pearse A. Keane<sup>13</sup>, Richard Savage<sup>42</sup>, Robert Golub<sup>43</sup>, Rupa Sarkar<sup>44</sup> and Samuel Rowley<sup>45</sup>**

<sup>14</sup>Department of Medicine, Women's College Research Institute, Women's College Hospital, University of Toronto, Toronto, Ontario, Canada. <sup>15</sup>Patient Safety Translational Research Centre, Imperial College London, London, UK. <sup>16</sup>Institute of Global Health Innovation, Imperial College London, London, UK. <sup>17</sup>Alan Turing Institute, London, UK. <sup>18</sup>Department of Statistics and Nuffield Department of Medicine, University of Oxford, Oxford, UK. <sup>19</sup>University of Manchester, Manchester, UK. <sup>20</sup>Department of Ophthalmology, Cantonal Hospital Lucerne, Lucerne, Switzerland. <sup>21</sup>University of Warwick, Coventry, UK. <sup>22</sup>Department of Ophthalmology, University of Washington, Seattle, WA, USA. <sup>23</sup>The National Institute for Health and Care Excellence, London, UK. <sup>24</sup>Salesforce Research, San Francisco, CA, USA. <sup>25</sup>Harvard T.H. Chan School of Public Health, Boston, MA, USA. <sup>26</sup>Medicines and Healthcare products Regulatory Agency, London, UK. <sup>27</sup>New England Journal of Medicine, Waltham, MA, USA. <sup>28</sup>Google Health, London, UK. <sup>29</sup>Annals of Internal Medicine, Philadelphia, PA, USA. <sup>30</sup>Patient Partner, Birmingham, UK. <sup>31</sup>British Medical Journal, London, UK. <sup>32</sup>National Institutes of Health, Bethesda, MD, USA. <sup>33</sup>Patient Partner, London, UK. <sup>34</sup>Patient Partner, Centre for Patient Reported Outcome Research, Institute of Applied Health Research, University of Birmingham, Birmingham, UK. <sup>35</sup>Centre for Statistics in Medicine, University of Oxford, Oxford, UK. <sup>36</sup>Hardian Health, London, UK. <sup>37</sup>AstraZeneca, Cambridge, UK. <sup>38</sup>Nature Research, New York, NY, USA. <sup>39</sup>Food and Drug Administration, Silver Spring, MD, USA. <sup>40</sup>Australian Institute for Machine Learning, North Terrace, Adelaide, Australia. <sup>41</sup>The Hospital for Sick Children, Toronto, Canada. <sup>42</sup>PinPoint Data Science, Leeds, UK. <sup>43</sup>Journal of the American Medical Association, Chicago, IL, USA. <sup>44</sup>The Lancet Group, London, UK. <sup>45</sup>Medical Research Council, London, UK.